

# Ensemble Methods (Kaggle winning entries)

Classifier 1 - KNN  
Classifier 2 - SUM  
Classifier 3 - DT

↑ KNN  
↑ SUM  
↑ DT

For an instance, predicts as: { Good, Good, Bad }

Majority voting  
Good.

KNN - 20  
SUM - 30  
DT - 30

} Ensemble will predict Good

when does ensemble learning not work well?

① Base model is bad

② All give similar or a prediction

Error of each model =  $l_y = .3$

$$\begin{aligned} P(2 \text{ model are wrong}) &= \binom{3}{2} (l_y)^2 (1-l_y)^{3-2} \\ &\quad + \binom{3}{3} (l_y)^3 (1-l_y)^{3-3} \\ &= .19 < .3 \end{aligned}$$

Single classifier (eg. decision tree)

Q: Feed in same data  $\rightarrow$  Different model / prediction?

Bagging (OR BOOTSTRAP AGGREGATION)

KEY IDEA: REDUCE VARIANCE  
But HOW??

Think about cross-validation<sup>2</sup>!

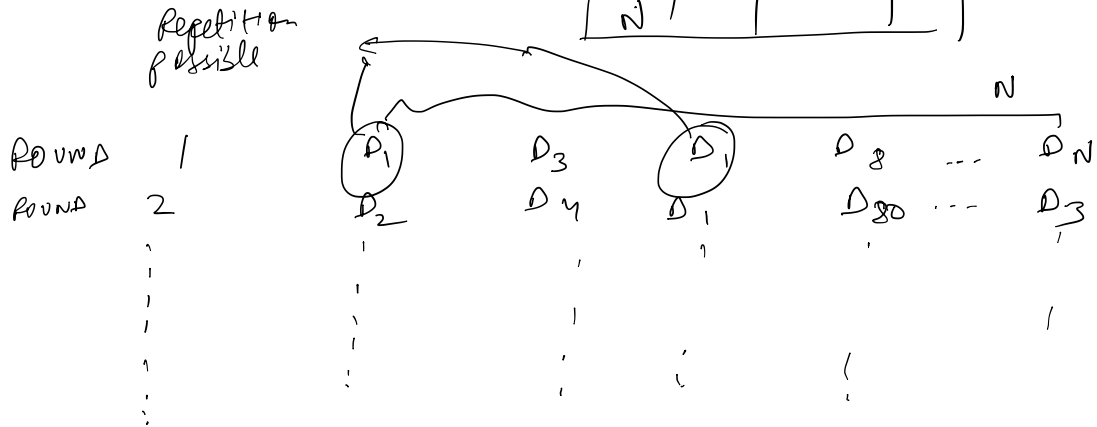
Similar in spirit  $\Rightarrow$  Create different datasets.

BUT, we have a single dataset!

SAMPLE WITH REPLACEMENT

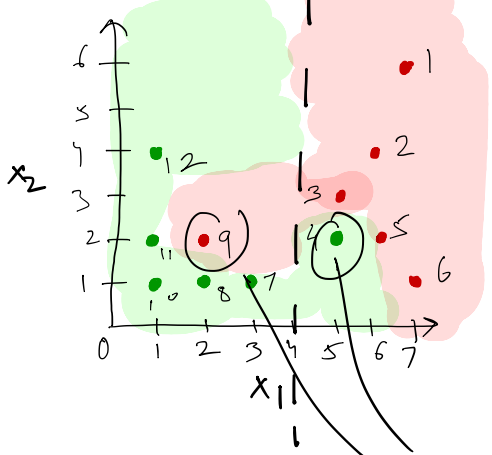
eg. if dataset is like.

Set	$F_1$	$F_2 \dots$	$y$
1			
2			
...			
N			

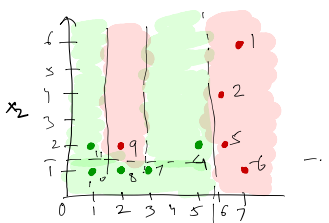


TRAIN SAME CLASSIFIER / REGRESSOR ON THESE DIFFERENT "DAGGING ROUNDS"

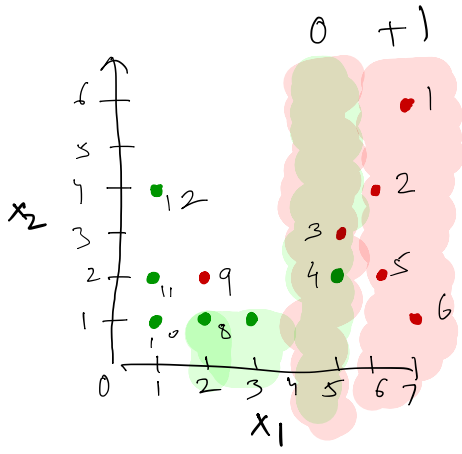
COMPLEX DT (HIGH VARIANCE)



(BAGGING ROUND #1  
DT 1 (3, 12 missing))

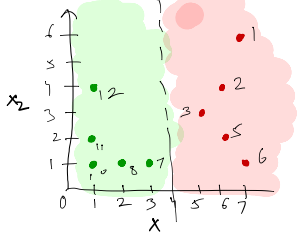


COMBINATION



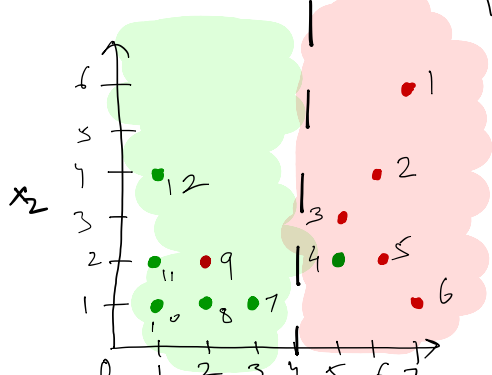
x lesser  
variance  
(even with  
full depth  
decision  
trees)

DT 2



BAGGING ROUND #2  
(9 and 12 missing)

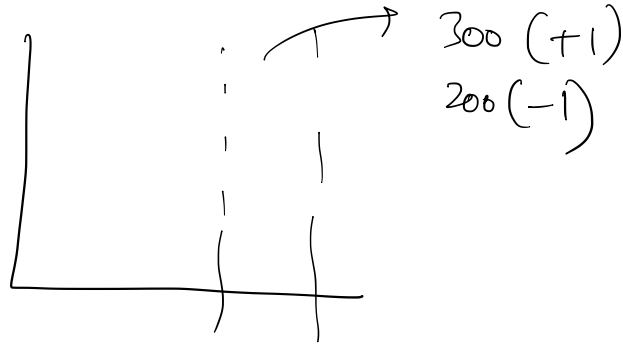
OPTIMAL DT (Reduce



variance  
by  
decreasing  
depth)

← we could  
learn if we  
limit depth = 1

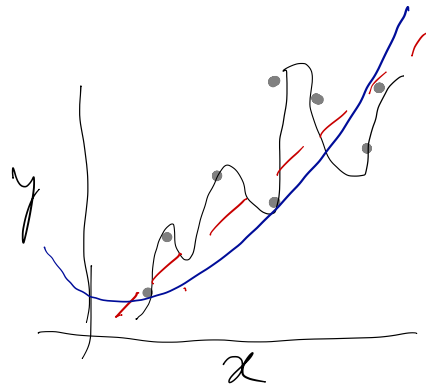
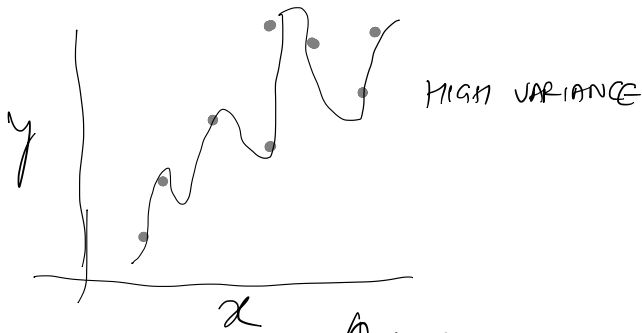
LABELS AT  
TEST  
TIME



$$\text{Signs } (300(+1) + 200(-1))$$

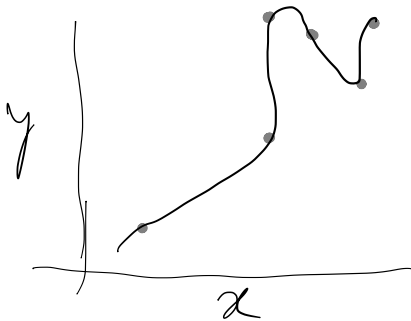
True  $\rightarrow$  <sup>Actual</sup> (+1)  
"RED"

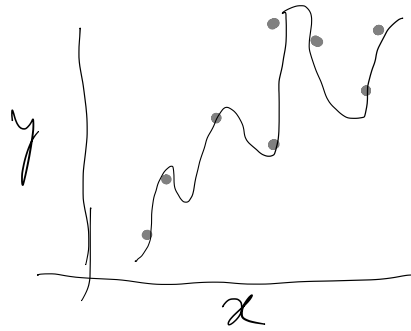
# BAGGING FOR REGRESSION



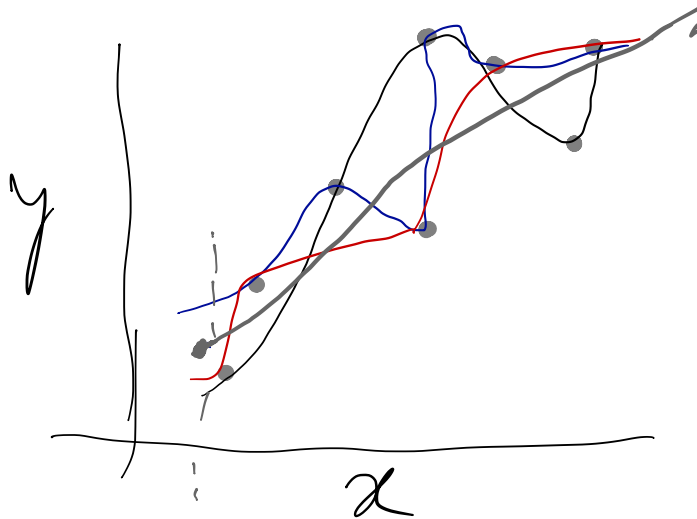
- Degra = 8
- Deg = 1
- Deg = 2

↑ Very different if some samples are removed..





ORIGINAL FIT



BAGGING  
Bagged ensemble



# Bagging

① Taking "strong" learners  
and combine  
(to reduce variance)

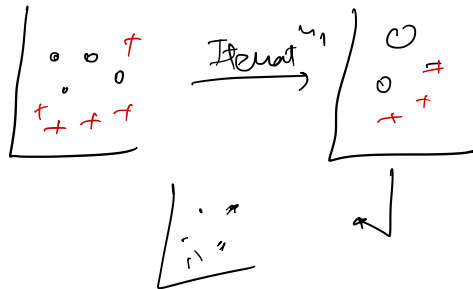
② Bagging: All learners are  
independent of each other...

"weigh samples  
differently"

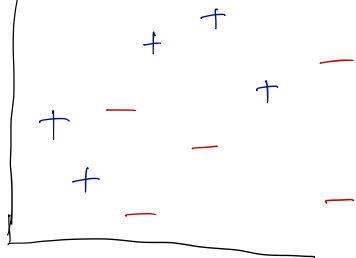
# Boosting

Taking "weak" learners  
(eg. depth-1  
 $\Delta \frac{1}{2}$ )  
and combine  
(to reduce  
bias)

Learners are  
incrementally built



# ADABOOST



TRAINING DATA

$$w_i \cdot i = \frac{1}{10} = .1$$

## ALGORITHM

① INIT.  $w_i^0 = \frac{1}{N}$

② FOR  $m=1$  TO  $M$ :

2.1) LEARN CLASSIFIER USING CURRENT WEIGHTS

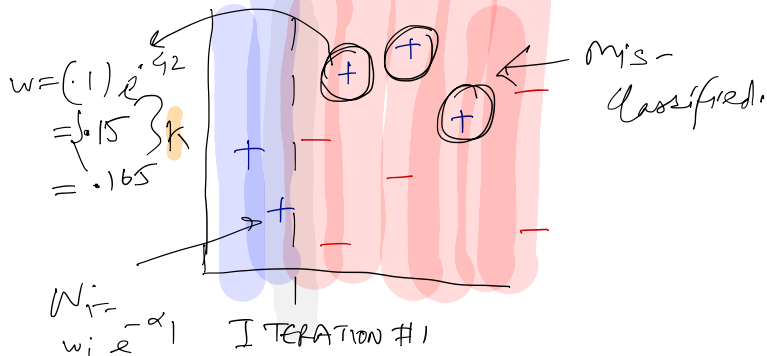
2.2) COMPUTE WEIGHTED ERROR =  $err_m = \frac{\sum w_i(\text{incorrect})}{\sum w_i}$

2.3) COMPUTE  $\alpha_m = \frac{1}{2} \log \left( \frac{1 - err_m}{err_m} \right)$

2.4) FOR CORRECT:  $w_i \leftarrow w_i e^{-\alpha_m}$

2.5) ' ' WRONG:  $w_i \leftarrow w_i e^{\alpha_m}$

2.6) NORMALIZE  $w_i$  TO SUM TO 1



$$w = (.1) e^{.42} = \{0.15\} * K = .071$$

$$err_1 = \frac{\sum w_i(\text{incorrectly})}{\sum w_i}$$

$$err_1 = \frac{3 \times .1}{10 \times .1} = .3$$

$$\alpha_1 = \frac{1}{2} \ln \left( \frac{1 - err_1}{err_1} \right) = \frac{1}{2} \ln \left( \frac{.7}{.3} \right) = \frac{1}{2} (.85) = .42$$

NORMALIZATION  $(K(.075) * 7 + K(.15) * 3) = 1$   
 $(.45 + .45) K = 1 \Rightarrow K = 1$

# ALGORITHM

① INIT.  $w_i^0 = \frac{1}{2}$

② FOR  $m=1$  TO  $M$ :

2.1) LEARN CLASSIFIER USING CURRENT WEIGHTS

2.2) COMPUTE WEIGHTED ERROR =  $err_m = \frac{\sum w_i(\text{incorrect})}{\sum w_i}$

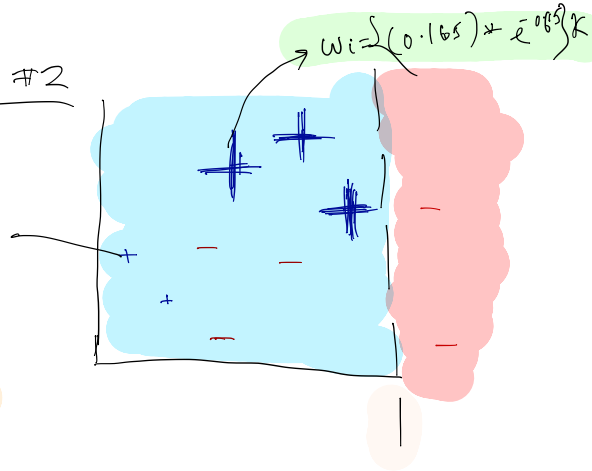
2.3) COMPUTE  $\alpha_m = \frac{1}{2} \log \left( \frac{1 - err_m}{err_m} \right)$

2.4) FOR CORRECT:  $w_i \leftarrow w_i e^{-\alpha_m}$

2.5) ' ' WRONG:  $w_i \leftarrow w_i e^{\alpha_m}$

2.6) NORMALIZE  $w_i$ s TO SUM TO 1

ITERATION #2

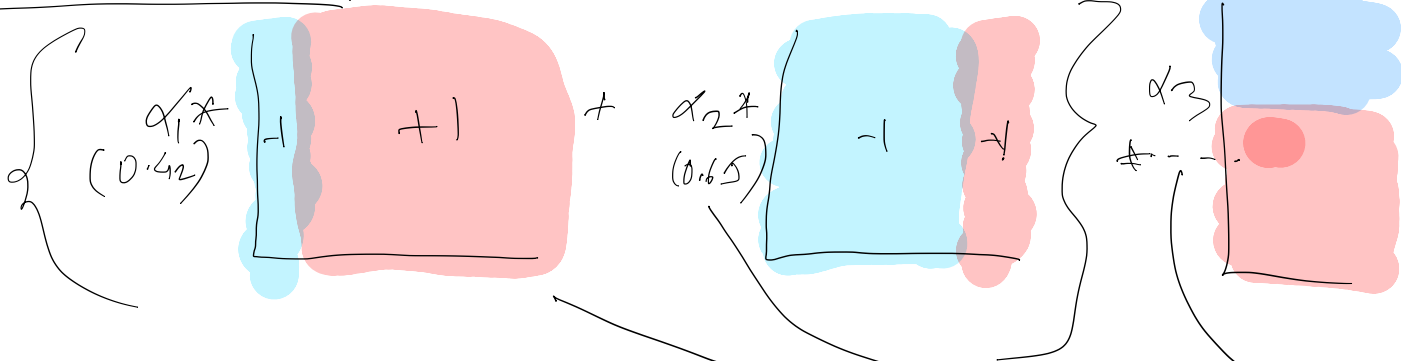


$dy_2 = 3 * 0.071 = 0.21$

$d_2 = \frac{1}{2} \log \left( \frac{.7}{.3} \right) = 0.65$

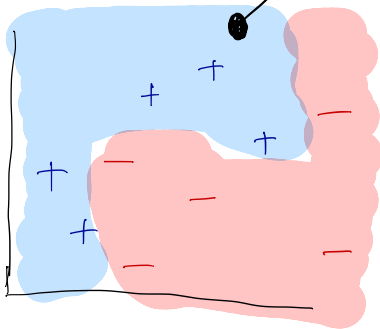
⋮

ADA BOOST TESTING

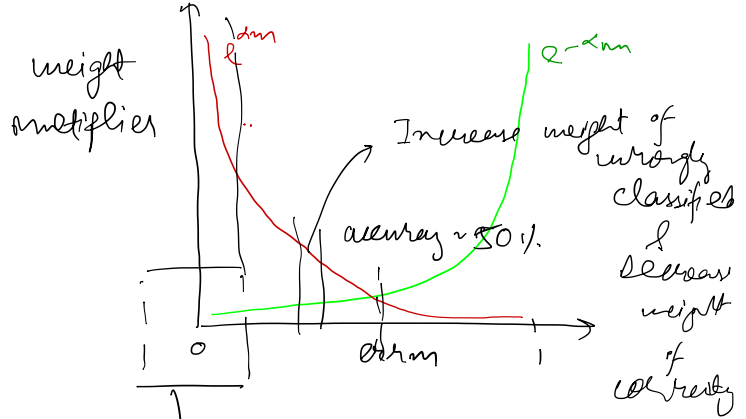
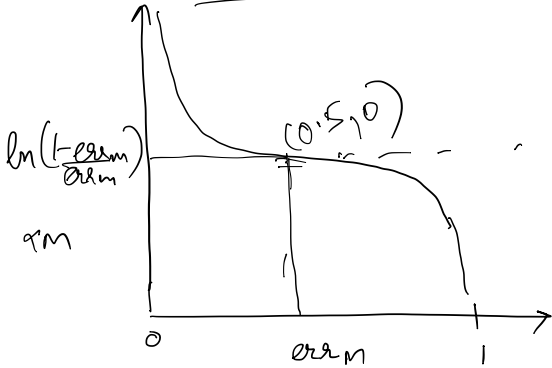


Effective classifier

$$0.42 \times (+1) + (0.65) \times (-1) + \frac{1}{3} \times (-1)$$



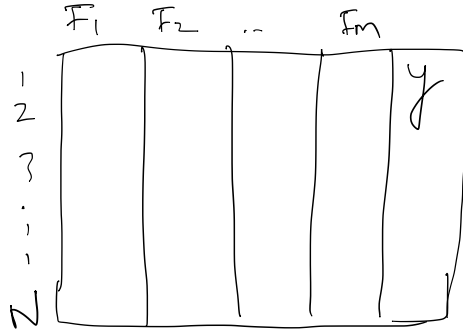
# WEIGHTS INTUITION



error  $\rightarrow 0$

Increase weights of- wrongly classified by a large amount.

# Random Forests



Reduce variance

↓  
↓  
Decorrelated trees

FOR ESTIMATOR  $m = 1$  to  $M$  - ESTIMATORS

① Create bootstrap dataset

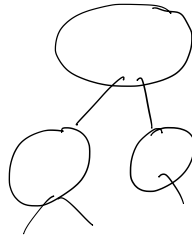
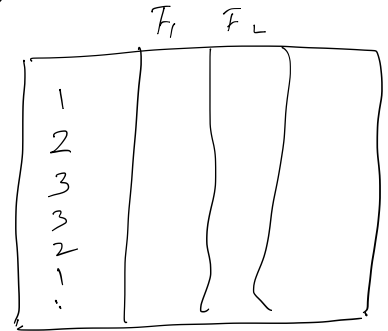
② WHILE GROWING TREE

②.1

②.2

CHOOSE  $m < M$  features for each node

Leads to



$F$ -AVAILABLE =  $[F_1, F_2, \dots]$

$F$ -AVAILABLE =  $[ \dots ]$



python script for

depth  
tree  
feature

① Depth

② # trees

③ #feature  $(\sqrt{M}, \sqrt{M+1}, \dots$   
 $\left. \begin{matrix} M_1, \dots, M_2, \dots, M \end{matrix} \right\}$

$< 1, 100, 1 \sum$   
 $< 1, 100, 2 \sum$   
?

for depth in [1, ...]:

for tree in [1, ... ~~M~~ trees]:

for feature in [-----]:

learn on train (depth, tree, feature)

